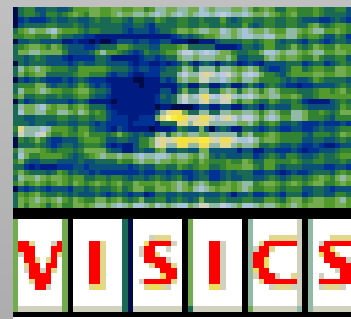# Pedestrian detection at 100 frames per second

R. Benenson, M. Mathias, R. Timofte and L. Van Gool

# Why 100 fps ?
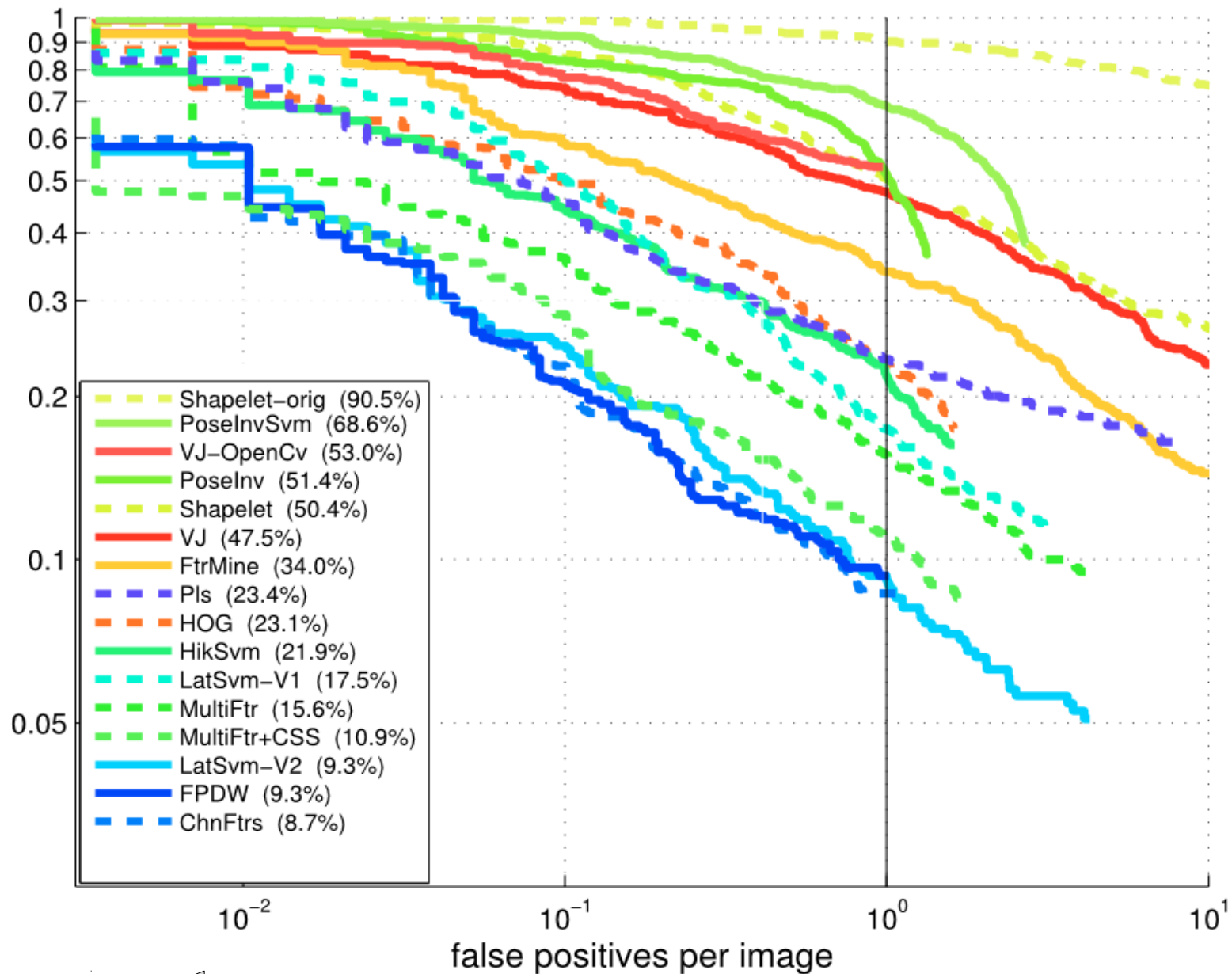
- Detection is one module amongst many

- Less computational power

- *Latency matters*

# How can we make things faster ?

# How can we make things *algorithmically* faster ?
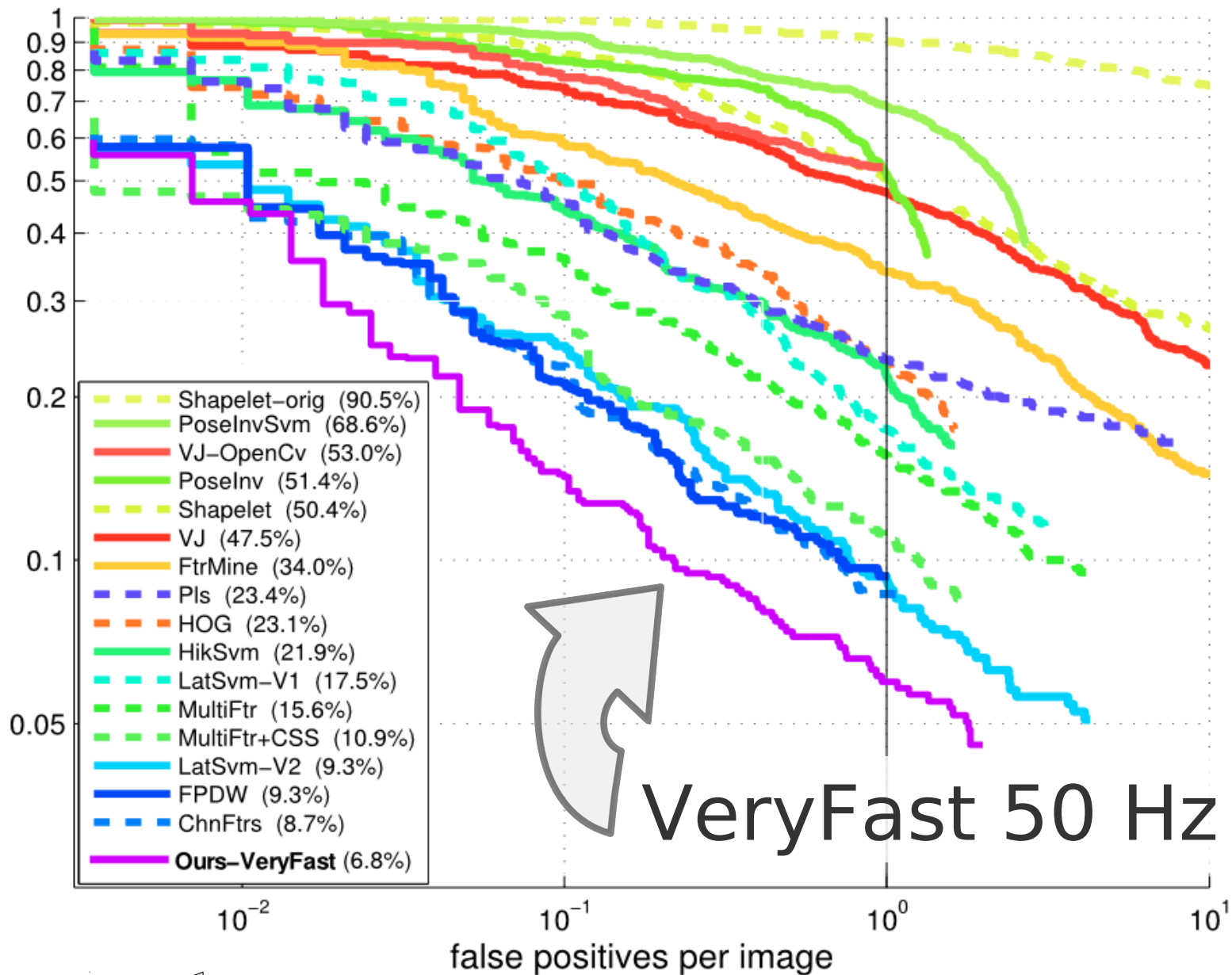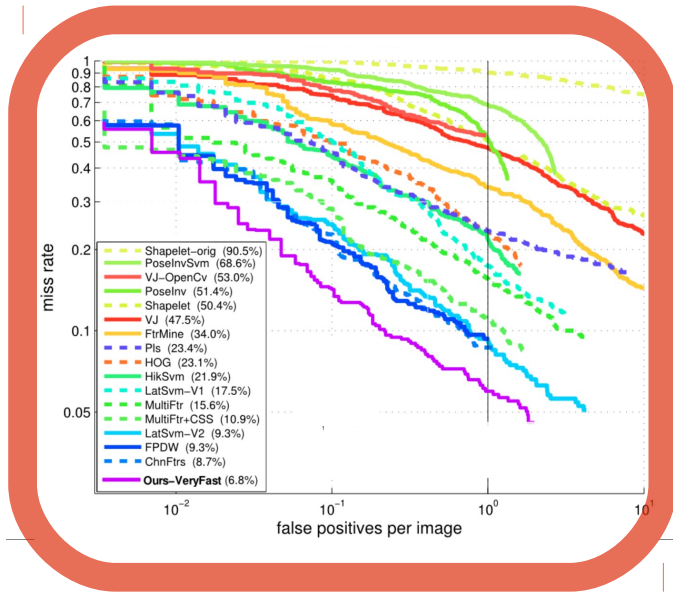
# INRIA dataset



miss rate

Better

- - - - Shapelet–orig (90.5%)
———— PoseInvSvm (68.6%)
———— VJ–OpenCv (53.0%)
———— PoseInv (51.4%)
- - - - Shapelet (50.4%)
———— VJ (47.5%)
———— FtrMine (34.0%)
- - - - Pls (23.4%)
- - - - HOG (23.1%)
———— HikSvm (21.9%)
- - - - LatSvm–V1 (17.5%)
- - - - MultiFtr (15.6%)
- - - - MultiFtr+CSS (10.9%)
———— LatSvm–V2 (9.3%)
———— FPDW (9.3%)
- - - - ChnFtrs (8.7%)

false positives per image

Better

[Dollar et al. 2011]

INRIA dataset

Better

Better

false positives per image

miss rate

- Shapelet−orig (90.5%)
- PoseInvSvm (68.6%)
- VJ−OpenCv (53.0%)
- PoseInv (51.4%)
- Shapelet (50.4%)
- VJ (47.5%)
- FtrMine (34.0%)
- Pls (23.4%)
- HOG (23.1%)
- HikSvm (21.9%)
- LatSvm−V1 (17.5%)
- MultiFtr (15.6%)
- MultiFtr+CSS (10.9%)
- LatSvm−V2 (9.3%)
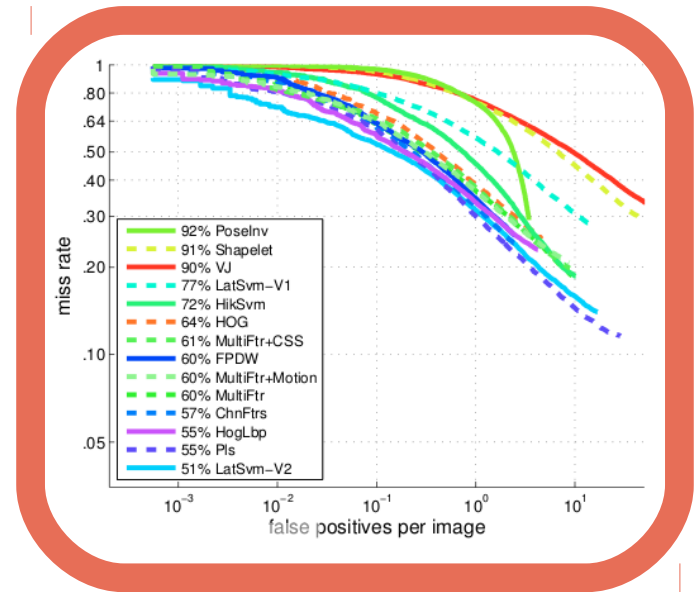- FPDW (9.3%)
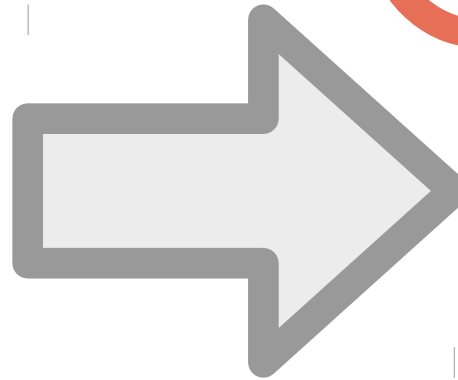- ChnFtrs (8.7%)
- **Ours−VeryFast (6.8%)**

VeryFast 50 Hz

INRIA dataset

ETH dataset

Monocular
50 Hz

Stereo
135 Hz

INRIA dataset

miss rate vs. false positives per image

- Shapelet−orig (90.5%)
- PoseInvSvm (68.6%)
- VJ−OpenCv (53.0%)
- PoseInv (51.4%)
- Shapelet (50.4%)
- VJ (47.5%)
- FtrMine (34.0%)
- Pls (23.4%)
- HOG (23.1%)
- HikSvm (21.9%)
- LatSvm−V1 (17.5%)
- MultiFtr (15.6%)
- MultiFtr+CSS (10.9%)
- LatSvm−V2 (9.3%)
- FPDW (9.3%)
- ChnFtrs (8.7%)

INRIA dataset

miss rate

false positives per image

Legend:
- Shapelet–orig (90.5%)
- PoseInvSvm (68.6%)
- VJ–OpenCv (53.0%)
- PoseInv (51.4%)
- Shapelet (50.4%)
- VJ (47.5%)
- FtrMine (34.0%)
- Pls (23.4%)
- HOG (23.1%)
- HikSvm (21.9%)
- LatSvm–V1 (17.5%)
- MultiFtr (15.6%)
- MultiFtr+CSS (10.9%)
- LatSvm–V2 (9.3%)
- FPDW (9.3%)
- ChnFtrs (8.7%)

[Viola & Jones 2004]

fastHOG
~10 Hz on GPU
[Prisacariu 2009]

Parts Model
[Felzenszwalb 2008]

ChnFtrs/FPDW
~5 Hz on CPU
[Dollar 2009+2010]

6 Orientation bins      Gradient magnitude      LUV color channels

6 Orientation bins

Gradient magnitude

LUV color channels

+1  -1  +1  -1

$$score = \quad w_1 \cdot h_1 +$$

6 Orientation bins

Gradient magnitude

LUV color channels

+1   -1   +1   -1          +1   -1   +1   -1
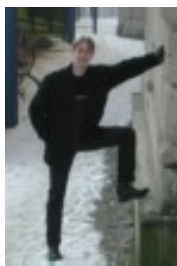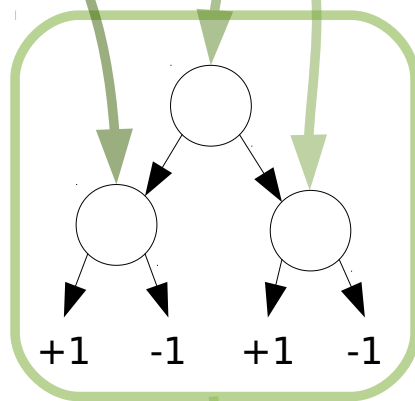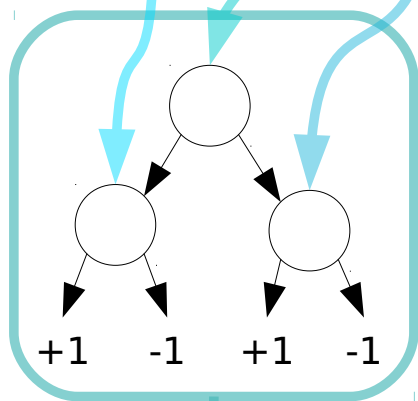
$$score = \quad w_1 \cdot h_1 + \qquad w_2 \cdot h_2 +$$

6 Orientation bins  Gradient magnitude  LUV color channels

$$score = \quad w_1 \cdot h_1 + \qquad w_2 \cdot h_2 + \qquad \cdots \qquad + w_N \cdot h_N$$

[ChnFtrs, Dollar et al. 2009]        (~4 Hz on GPU)

# What slows down fastHOG ?



GPU Time (Total, %)

- Block histogram computation
- Linear SVM evaluation
- Gradient magnitudes and orientations
- Block histogram normalization
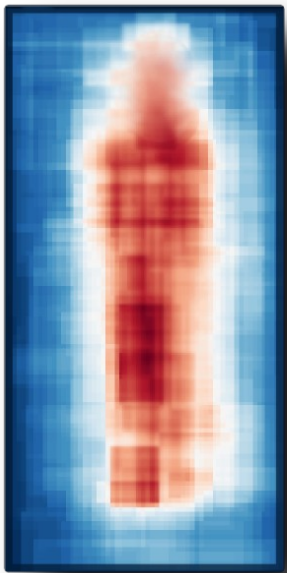- Image downscale
- Others
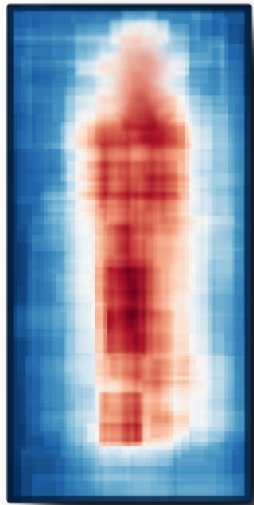
3.14%　0.76%
3.44%
11.69%
16.22%
64.75%

[Prisacariu and Reid 2009]

# How to make *features computation* faster ?

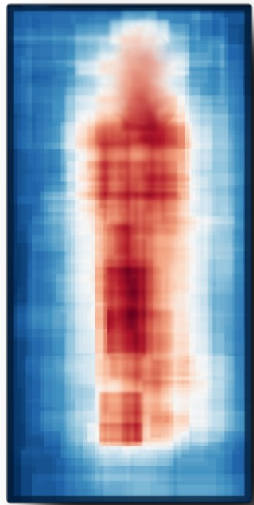# One template cannot detect at multiple scales

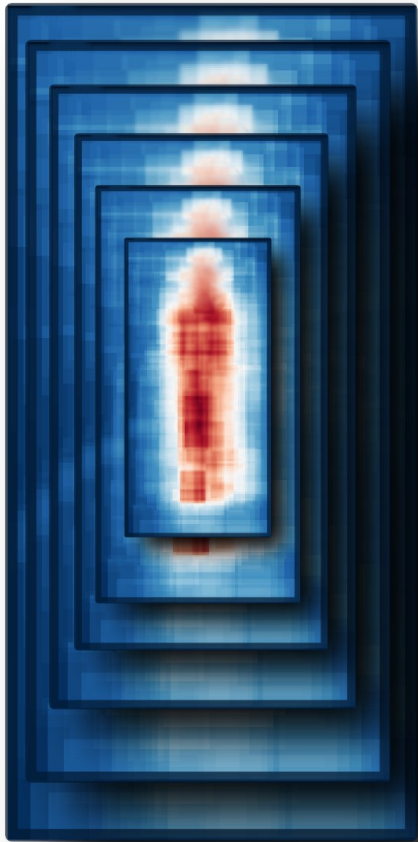# Traditionally, features are computed many times



~50 scales

# Traditionally, features are computed many times



~50 scales

# Training one model per scale is too expensive
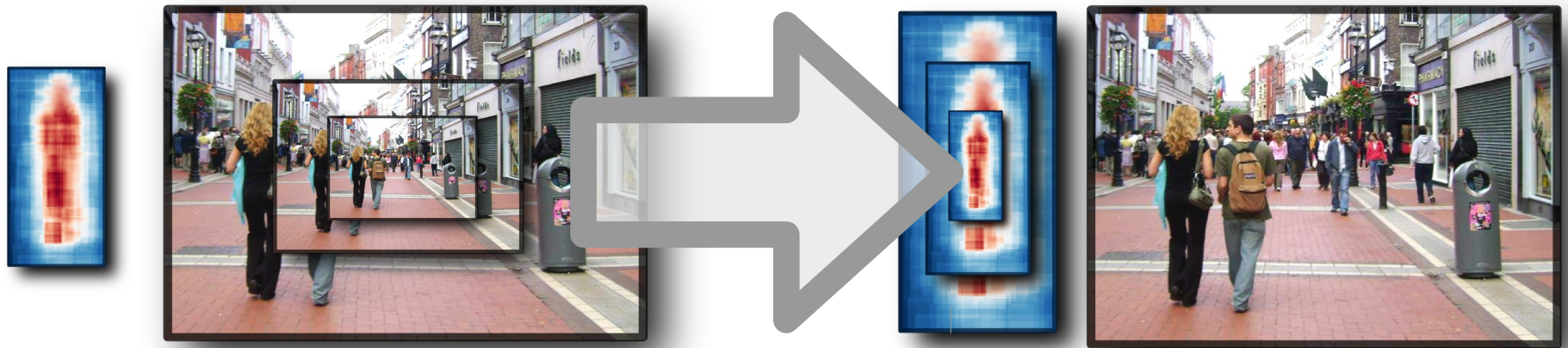


~50 scales

# Can we avoid resizing
the input image many times ?

# Features can be approximated across scales



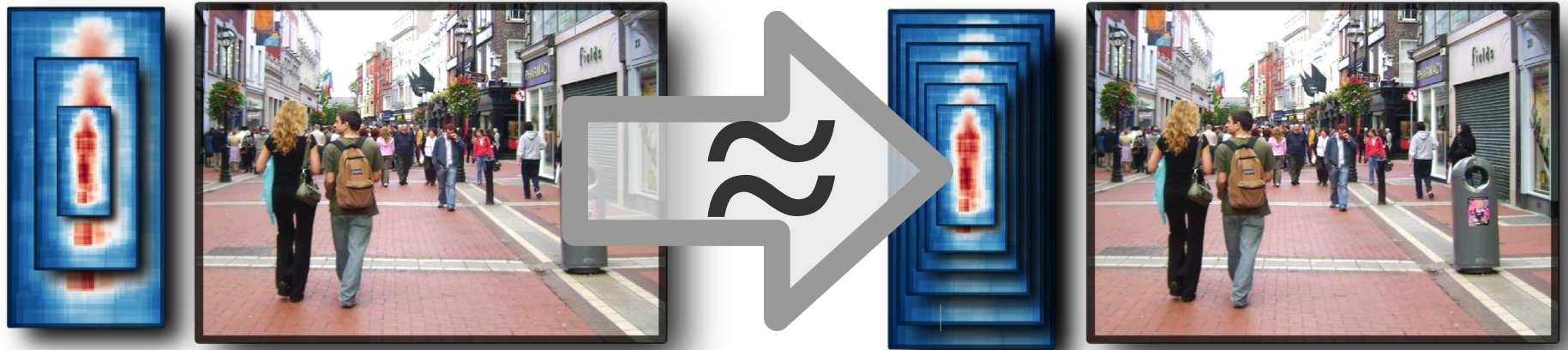~5 scales          ~50 scales

[Dollar et al. 2010]

# We transfer test time computation to training time
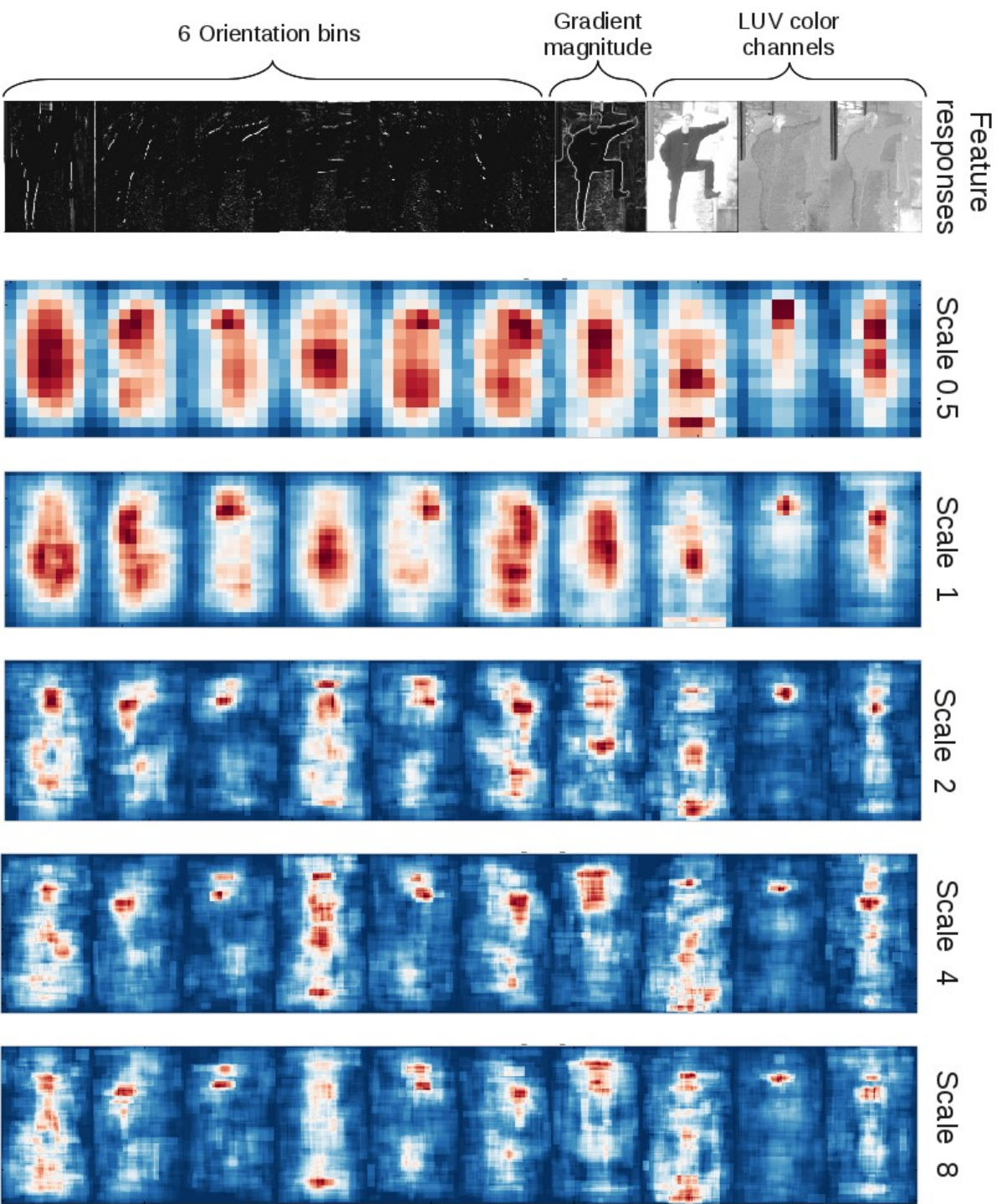


1 model,
5 image scales

5 models,
**1 image scale**

(3x reduction in features computation)

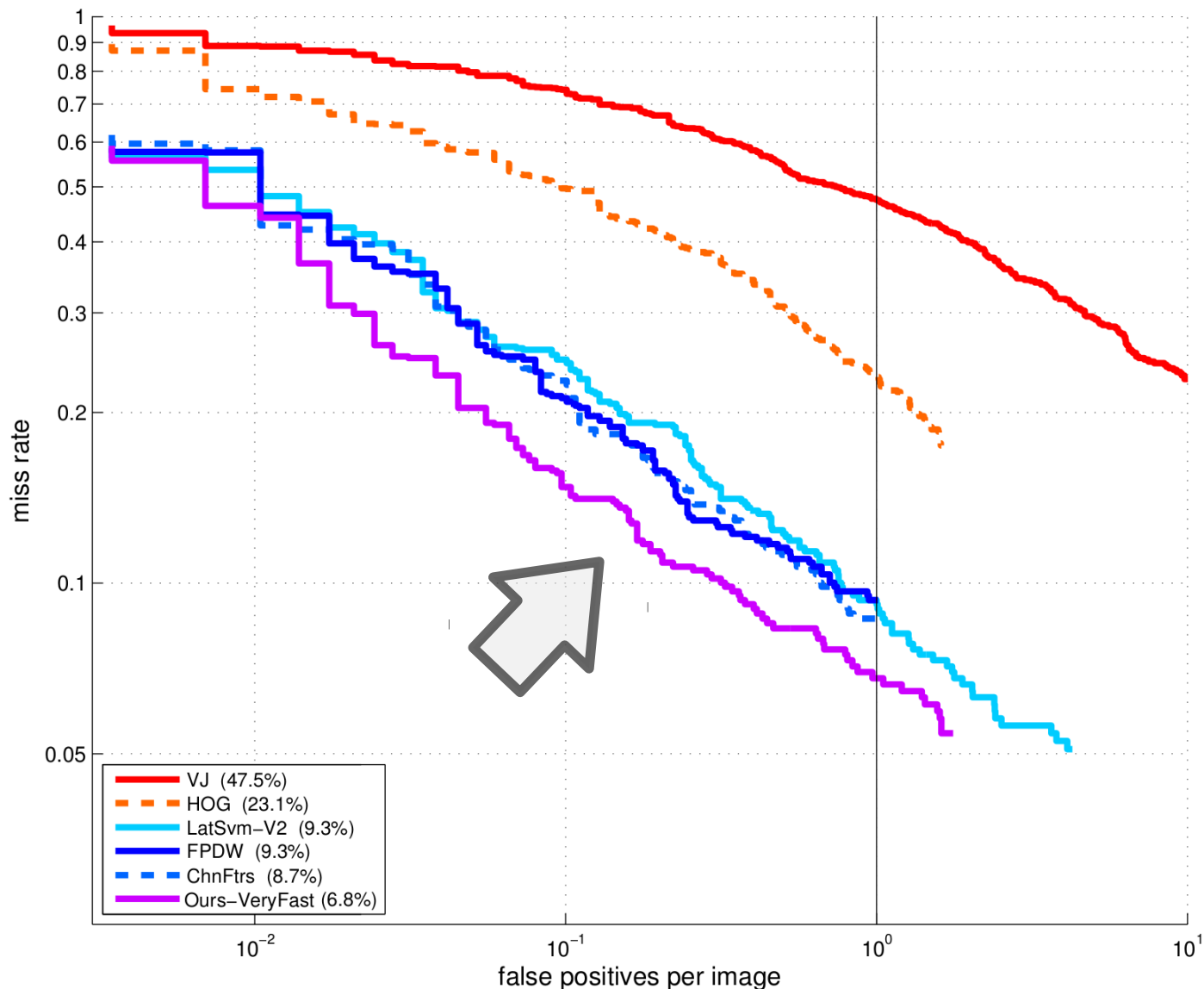# At runtime, we use as many models as scales



5 models,
1 image scale

50 models,
1 image scale

6 Orientation bins · Gradient magnitude · LUV color channels

Feature responses

Scale 0.5

Scale 1

Scale 2

Scale 4

Scale 8

# Detecting without resizing provides quality



miss rate

false positives per image

VJ (47.5%)
HOG (23.1%)
LatSvm−V2 (9.3%)
FPDW (9.3%)
ChnFtrs (8.7%)
Ours−VeryFast (6.8%)

# Detecting without resizing provides speed

- ~3x less time on features computation

- Avoids alternating between features and detection scores computation (relevant in practice)
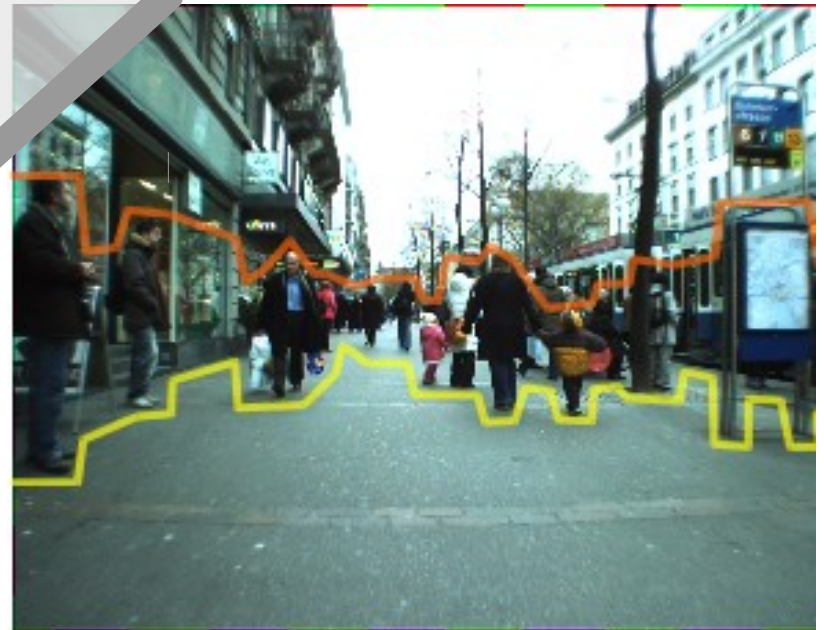
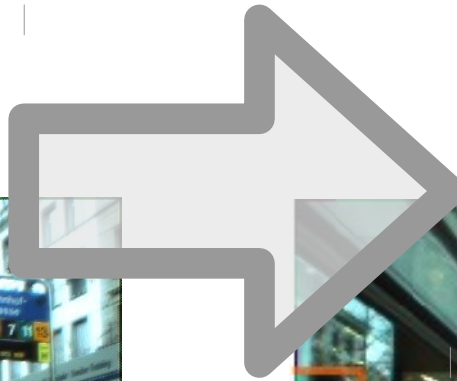- We reach **50 Hz** on GPU, 640x480 pixels x 55 scales

# We want to use scene geometry to guide the detections


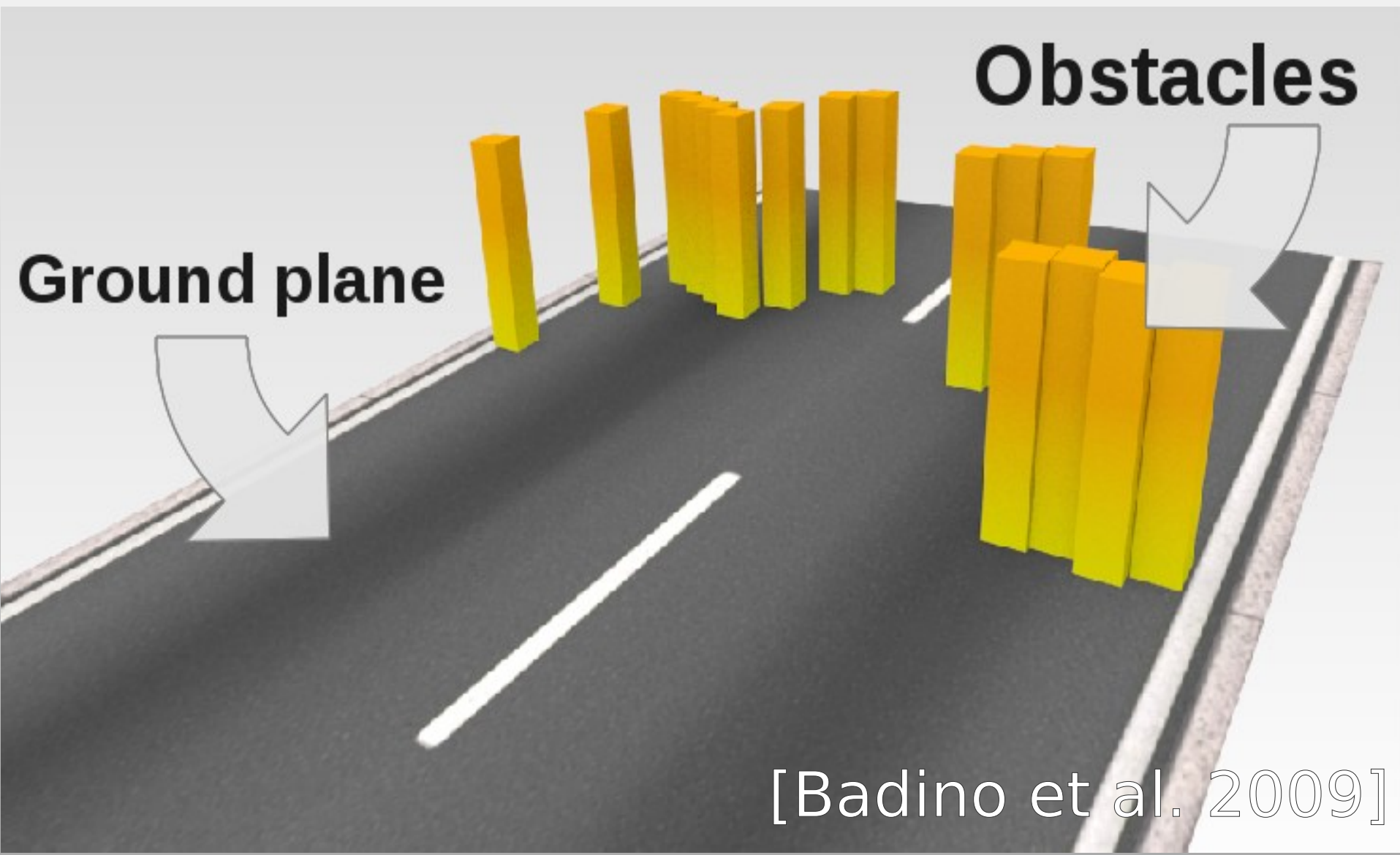
Monocular

Stereo

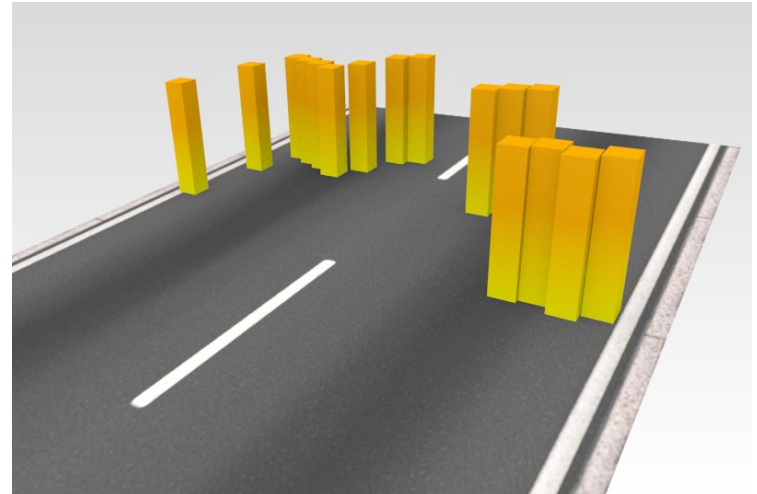# Stixel world
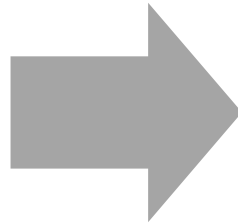


Obstacles

Ground plane

[Badino et al. 2009]

# Depth maps are slow to compute



<50 Hz on CPU

# Stixel world without depth map



**135 Hz** on CPU

[Benenson et al. 2011]

# Stixel world without depth map

# Stixel world without depth map



Left

Right

v

disparity

## 1) Ground plane estimation

# Stixel world without depth map



Left

Right

v

disparity

1) Ground plane estimation

# Stixel world without depth map



Left

Right

disparity

u

2) Stixel distance estimation

# Stixel world without depth map



Left

Right

disparity

u

2) Stixel distance estimation

# Stixel world without depth map



Left

Right

disparity

u

2) Stixel distance estimation @ 135 Hz CPU

# ETH's dataset results have less variance than INRIA's



Monocular

Stereo

# Using stixels provides speed without quality loss



Recall versus FPPI over Bahnhof dataset, considering all windows with height > 40 [pixels]

Better

Recall

False positives per image

- - - HOG + SVM Dalal and Triggs
—— VeryFast
- · - VeryFast + ground plane
—— VeryFast + stixels

# Detecting using stixels provides speed

- No geometry:
  - 640x480 pixels x 55 scales => 50 Hz on desktop

- Ground plane:
  - 640x60 pixels x 55 scales => 100 Hz on desktop (8x reduction in search space)

- Stixels:
  - 640x60 pixels x 10 scales => **135 Hz** on desktop (150 Hz GPU side, 135 Hz CPU side)
  - 44x reduction in search space
  - We reach *80 Hz on laptop*

# Win-win detector

- **High**est known **quality** for a single part detector
  (over the INRIA dataset, at camera ready time)

- 50 Hz in monocular mode,
  **135 Hz** in **stereo mode**, *80 Hz on a laptop*.

- 5x faster and 3x lower missrate
  than previous state-of-the-art, fastHOG.

# No resizing + stixels
# ==
# faster and better detections

# Future work

- Transforming classifier seems useful:
  - Extension for different occlusions (submitted)
  - Extension for different point of views ?

# Future work

- Transforming classifier seems useful:
    - Extension for different occlusions (submitted)
    - Extension for different point of views ?


- There is room for speed and quality improvements
    - (Original implementation was crude,
      work in progress version reaches ~170 Hz)

# Future work

- Transforming classifier seems useful:
  - Extension for different occlusions (submitted)
  - Extension for different point of views ?

- There is room for speed and quality improvements
  - (Original implementation was crude, work in progress version reaches ~170 Hz)

- Building a part pased model
  on top of our VeryFast detector ?

# ?

# Rodrigo Benenson
http://rodrigob.github.com

# Source code release on August 1st

**CPU**

**GPU**

$t_0$ Input images → Stereo rectification

Stereo rectification → Stixels computation

Stereo rectification → Pedestrian detection

Stixels computation → Pedestrian detection

Pedestrian detection → Non-maxima suppression

Non-maxima suppression → Result detections

7.4 ms
*(1000 frames average)*

$t_1$ Input images → Stereo rectification

Stereo rectification → Stixels computation

Stereo rectification → Pedestrian detection

Stixels computation → Pedestrian detection

6 cores of Intel Core i7-2630QM @ 2.00GHz

Nvidia GeForce GTX 560M